

JC777 U.S. PTO  
03/31/00

04-03-00

A

# IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Inventorship.....Brill et al.  
Applicant.....Microsoft Corporation  
Attorney's Docket No. ....MS1-471US  
Title: Spell Checker with Arbitrary Length String-to-String Transformations to Improve Noisy Channel  
Spelling Correction

JC678 U.S. PTO  
09/539357  
03/31/00

## TRANSMITTAL LETTER AND CERTIFICATE OF MAILING

To: Commissioner of Patents and Trademarks,  
Washington, D.C. 20231

From: Lewis C. Lee (Tel. 509-324-9256; Fax 509-323-8979)  
Lee & Hayes, PLLC  
421 W. Riverside Avenue, Suite 500  
Spokane, WA 99201

The following enumerated items accompany this transmittal letter and are being submitted for the matter identified in the above caption.

1. Specification—title page, plus 32 pages, including 57 claims and Abstract
2. Transmittal letter including Certificate of Express Mailing
3. 5 Sheets Formal Drawings (Figs. 1-6)
4. Return Post Card

Large Entity Status ☒ [x]

Small Entity Status ☐ [ ]

Date: March 31, 2000

By: 

Lewis C. Lee  
Reg. No. 34,656

## CERTIFICATE OF MAILING

I hereby certify that the items listed above as enclosed are being deposited with the U.S. Postal Service as either first class mail, or Express Mail if the blank for Express Mail No. is completed below, in an envelope addressed to The Commissioner of Patents and Trademarks, Washington, D.C. 20231, on the below-indicated date. Any Express Mail No. has also been marked on the listed items.

**EL580803682**

Express Mail No. (if applicable) \_\_\_\_\_

Date: March 31, 2000

By: 

Helen M. Hare

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

**Spell Checker with Arbitrary Length String-to-String  
Transformations to Improve Noisy Channel Spelling  
Correction**

Inventor(s):  
Eric Brill  
Robert Moore

ATTORNEY'S DOCKET NO. MS1-471US

## TECHNICAL FIELD

This invention relates to spell checkers used in computer programs to identify and potentially correct misspelled words.

## BACKGROUND

Spell checkers are well-known program components used in computer programs to inform users that a word is misspelled, and in some cases, to correct the error to the appropriate spelling. Word processing programs, email programs, spreadsheets, browsers, and the like are examples of computer programs that employ spell checkers.

One conventional type of spell checker corrects errors in an ad-hoc fashion by manually specifying the types of allowable edits and the weights associated with each edit type. For the spell checker to recognize an entry error "fysical" and correct the error to the appropriate word "physical", a designer manually specifies a substitution edit type that allows substitution of the letters "ph" for the letter "f". Since it is built manually, this approach does not readily port to a new language or adapt to an individual's typing style.

Another type of spell checker is one that learns errors and weights automatically, rather than being manually configured. One type of trainable spell checker is based on a noisy channel model, which observes character strings actually entered by a user and attempts to determine the intended string based on a model of generation.

Spell checkers based on the noisy channel model have two components: (1) a word or source generation model, and (2) a channel or error model. The source model describes how likely a particular word is to have been generated. The error

1 model describes how likely a person intending to input X will instead input Y.  
2 Together, the spell checker attempts to describe how likely a particular word is to  
3 be the intended word, given an observed string that was entered.

4 As an example, suppose a user intends to type the word "physical", but  
5 instead types "fysical". The source model evaluates how likely the user is to have  
6 intended the word "physical". The error model evaluates how likely the user is to  
7 type in the erroneous word "fysical" when the intended word is "physical".

8 The classic error model computes the Levenshtein Distance between two  
9 strings, which is the minimum number of single letter insertions, deletions, and  
10 substitutions needed to transform one character string into another. The classic  
11 error model is described in Levenshtein, V. "Binary Codes Capable of Correcting  
12 Deletions, Insertions and Reversals." Soviet Physics – Doklady 10, 10, pp. 707-  
13 710. 1966.

14 A modification of the classic error model employs a Weighted Levenshtein  
15 Distance, in which each edit operation is assigned a different weight. For instance,  
16 the weight assigned to the operation "Substitute e for i" is significantly different  
17 than the weight assigned to the operation "Substitute e for M". Essentially all  
18 existing spell checkers that are based on edit operations use the weighted  
19 Levenshtein Distance as the error model, while sometimes adding a small number  
20 of additional edit templates, such as transposition, doubling, and halving.

21 The error model can be implemented in several ways. One way is to  
22 assume all edits are equally likely. In an article by Mays, E., Damerau, F, and  
23 Mercer, R. entitled "Context Based Spelling Correction," Information Processing  
24 and Management, Vol. 27, No. 5, pp. 517-522, 1991, the authors describe pre-  
25 computing a set of edit-neighbors for every word in the dictionary. A word is an

1 edit-neighbor of another word, if it can be derived from the other word from a  
2 single edit, where an edit is defined as a single letter insertion (e.g.,  $\emptyset \rightarrow a$ ), a  
3 single letter substitution (e.g.,  $a \rightarrow b$ ), a single letter deletion (e.g.,  $a \rightarrow \emptyset$ ), or a  
4 letter-pair transposition (e.g.,  $ab \rightarrow ba$ ). For every word in a document, the spell  
5 checker determines whether any edit-neighbor of that word is more likely to  
6 appear in that context than the word that was typed. All edit-neighbors of a word  
7 are assigned equal probability of having been the intended word, and the context is  
8 used to determine which word to select. It is noted that the word itself (if it is in  
9 the dictionary) is considered an edit-neighbor of itself, and it is given a much  
10 higher probability of being the intended word than the other edit-neighbors.

11 A second way to implement the error model is to estimate the probabilities  
12 of various edits from training data. In an article by Church, K. and Gale, W.,  
13 entitled "Probability Scoring for Spelling Correction," Statistics and Computing  
14 1, pp. 93-103, 1991, the authors propose employing the identical set of edit types  
15 used by Mays et al. (i.e., single letter insertion, substitution, deletion, and letter-  
16 pair transposition) and automatically deriving probabilities for all edits by  
17 computing the probability of an intended word  $w$  given an entered string  $s$ . The  
18 Church et al. method trains on a training corpus to learn the probabilities for each  
19 possible change, regardless of the correct word and entered word. In other words,  
20 it learns the probability that an erroneous input string  $s$  will be written when the  
21 correct word  $w$  was intended, or  $P(s|w)$ . The Church et al. method improves  
22 insertion and deletion by including one character of context.

23 The error model probability  $P(s|w)$  used in noisy channel spell correction  
24 programs, such as the one described in Church et al., may seem backwards  
25 initially because it suggests finding how likely a string  $s$  is to be entered given that

1 a dictionary word  $w$  is intended. In contrast, the spell correction program actually  
2 wants to know how likely the entered string  $s$  is to be a word  $w$  in the dictionary,  
3 or  $P(w|s)$ . The error model probability  $P(s|w)$  comes from Bayes formula, which  
4 can be used to represent the desired probability  $P(w|s)$  as follows:

$$P(w|s) = \frac{P(s|w) \cdot P(w)}{P(s)}$$

8 The denominator  $P(s)$  remains the same for purposes of comparing possible  
9 intended words given the entered string. Accordingly, the spell checking analysis  
10 concerns only the numerator product  $P(s|w) \cdot P(w)$ , where the probability  $P(s|w)$   
11 represents the error model and the probability  $P(w)$  represents the source model.

12 As application programs become more sophisticated and the needs of users  
13 evolve, there is an ongoing need to improve spell checkers. The inventors have  
14 developed an improved spell checker that is based on the noisy channel model,  
15 which incorporates a more powerful error model component.

## 17 SUMMARY

18 A spell checker based on the noisy channel model has a source model and  
19 an error model. The source model determines how likely a word  $w$  in a dictionary  
20 is to have been generated. The error model determines how likely the word  $w$  was  
21 to have been incorrectly entered as the string  $s$  (e.g., mistyped or incorrectly  
22 interpreted by a speech recognition system).

23 The error model determines this probability based on edit operations that  
24 convert arbitrary length character sequences in the word  $w$  to arbitrary length  
25

1 character sequences in the string  $s$ . These edit operations are characterized as  $\alpha \rightarrow$   
2  $\beta$ , where  $\alpha$  is one character sequence of zero or more characters and  $\beta$  is another  
3 character sequence of zero or more characters. In many cases, the number of  
4 characters in each sequence  $\alpha$  and  $\beta$  will be different. In this manner, the edit  
5 operations are not constrained or limited to the specified set of changes, such as  
6 single letter insertion, deletion, or substitution.

7 The error model determines how likely a word  $w$  in the dictionary was to  
8 have been mistyped as the string  $s$  (i.e.  $P(s|w)$ ) according to the probabilities of the  
9 string-to-string edits. One implementation is to find all possible sets of string-to-  
10 string edits that transform the word  $w$  into the string  $s$ , calculating  $P(s|w)$  for each  
11 set and then summing over all sets. The probabilities are derived through a  
12 training process that initially uses Levenshtein Distance or other cost metric to  
13 find the least cost alignment of characters in a pair of wrong and right inputs.

## 14 **BRIEF DESCRIPTION OF THE DRAWINGS**

15  
16 Fig. 1 is a block diagram of an exemplary computer that runs a spell  
17 checker.

18 Fig. 2 is a flow diagram of a process implemented by the spell checker to  
19 compute a probability  $P(s|w)$  given an erroneously entered string  $s$  and a dictionary  
20 word  $w$ .

21 Fig. 3 is a block diagram of a training computer used to train the spell  
22 checker.

23 Fig. 4 is a flow diagram of a training method implemented by the training  
24 computer.

1 Fig. 5 is a diagrammatic illustration of an alignment technique used during  
2 training of an error model employed in the spell checker.

3 Fig. 6 is a diagrammatic illustration of the alignment technique of Fig. 5 at  
4 a point later in the process.

## 5 6 **DETAILED DESCRIPTION**

7 This invention concerns a spell checker used in computer programs to  
8 identify and, in some cases, correct misspelled words. The spell checker may be  
9 used in many different applications, including word processing programs, email  
10 programs, spreadsheets, browsers, and the like. For discussion purposes, the spell  
11 checker is described in the context of a spell correction program implemented in a  
12 word processing program.

13 However, aspects of this invention may be implemented in other  
14 environments and in other types of programs. For instance, the invention may be  
15 implemented in language conversion software (e.g., Japanese Katakana to English)  
16 that may implement string-to-string mapping.

## 17 18 **Exemplary Computing Environment**

19 Fig. 1 shows an exemplary computer 20 having a processor 22, volatile  
20 memory 24 (e.g., RAM), and non-volatile memory 26 (e.g., ROM, Flash, hard  
21 disk, floppy disk, CD-ROM, etc.). The computer 20 also has one or more input  
22 devices 28 (e.g., keyboard, mouse, microphone, stylus, etc.) and a display 30 (e.g.,  
23 monitor, LCD, etc.). The computer may also have other output devices (now not  
24 shown), such as a speaker. The computer 20 is representative of many diverse  
25



1 types of computing devices, including desktop computers, laptops, handheld  
2 computers, set-top boxes, information appliances, and so forth.

3 The computer 20 runs an operating system 32 and a word processing  
4 application program 34. For purposes of illustration, operating system 32 and  
5 word processor 34 are illustrated herein as discrete blocks stored in the non-  
6 volatile memory 26, although it is recognized that such programs and components  
7 reside at various times in different storage components of the computer 20 and are  
8 executed by the processor 22. Generally, these software components are stored in  
9 non-volatile memory 26 and from there, are loaded at least partially into the  
10 volatile main memory 24 for execution on the processor 22.

11 The word processor 34 includes a spell correction program 40 that  
12 identifies and, where appropriate, corrects misspelled words that the user has  
13 entered. The user enters the words in a conventional manner, such as typing them  
14 in on a keyboard, using a stylus to input individual letters, or speaking words into  
15 a microphone. In the event voice entry is employed, the computer 20 would also  
16 implement a speech recognition program (not shown) to convert voice input to  
17 words.

18 The spell correction program 40 is based on the noisy channel model and  
19 has two components: a source model 42 and an error model 44. The source model  
20 42 includes computer-executable instructions that determine how likely a  
21 particular word  $w$  in a dictionary  $D$  is to have been generated in this particular  
22 context. The source model is represented statistically as the probability  $P(w |$   
23  $context)$ . The error model 44 includes computer-executable instructions that  
24 determine how likely a user is to enter the string  $s$  when intending to enter the  
25 word  $w$ . The error model is represented statistically as the probability  $P(s|w)$ .

1 Accordingly, the spell checker 40 attempts to correct an erroneously  
2 entered string  $s$  into a word  $w$  by returning the change that maximizes the  
3 probabilities of the source and error models, as follows:

$$\arg \max_{w \in D} P(s | w) \times P(w | context)$$

4  
5  
6  
7 The source and error models are independent of each other. Thus, different  
8 source models 42 may be used interchangeably with the preferred error model  
9 described below. The source model 42 may be implemented as a language model,  
10 if one is available. In this case,  $P(w | context)$  is the probability of word  $w$ , given  
11 the context words that the user had entered prior to  $w$ . If the relative word  
12 probabilities are known, but not anything about how the words are used in context,  
13  $P(w | context)$  can be reduced to  $P(w)$ . Finally, if the spell checker knows nothing  
14 about the relative word probabilities of the user generating words in the dictionary,  
15 the source model can be eliminated entirely by setting  $P(w | context)$  to  $\frac{1}{|D|}$  for all  
16  $w$ .

#### 17 18 **Improved Error Model 44**

19 Unlike conventional error models, the error model 44 employed in spell  
20 checker 40 permits edit operations that convert a first string of arbitrary size to a  
21 second string of arbitrary size. That is, given an alphabet  $V$  and arbitrary length  
22 character sequences  $\alpha$  and  $\beta$ , the error model 44 allows edit operations of the form  
23  $\alpha \rightarrow \beta$ , where  $\alpha, \beta \in V^*$  (where  $V^*$  represents the set of all strings of characters in  
24  $V$  of length 0 or more). Each character sequence  $\alpha$  and  $\beta$  may have zero or more  
25 characters and in many cases, the number of characters in each sequence  $\alpha$  and  $\beta$

1 will be different. In this manner, the edit operations are not constrained or limited  
2 to a specified set of changes, such as single letter insertion, deletion, or  
3 substitution.

4 The error model 44 can therefore be characterized as the probability that,  
5 when a user intends to type a character sequence  $\alpha$ , he/she instead types  $\beta$ . This is  
6 represented as the probability  $P(\beta|\alpha)$ . The error model probability  $P(s|w)$  can then  
7 be expressed as a set of probabilities describing various arbitrary length string-to-  
8 string conversions, as follows:

$$P(s|w) = P(\beta_1|\alpha_1) * P(\beta_2|\alpha_2) * P(\beta_3|\alpha_3) * \dots * P(\beta_n|\alpha_n)$$

12 One implementation of the error model 44 is to find all possible sets of  
13 string-to-string edits that transform the word  $w$  into the string  $s$ , calculate the  
14 probability  $P(s|w)$  for each set using the above formula, and sum over all sets.  
15 More particularly, for each possible word  $w$  that the erroneous string  $s$  might be,  
16 the error model 44 partitions the word  $w$  and string  $s$  into different numbers of  
17 segments that define varying lengths of character sequences. For example,  
18 suppose the dictionary word is "physical" and the number of partition segments is  
19 five. One possible partition is, say, "ph y s ic al". Now, suppose that a user  
20 generates each partition, possibly with errors. One possible result of the user input  
21 is a string "fisikle" with a five-segment partition "f i s ik le".

22 The error model then computes the probability  $P(\beta|\alpha)$  for each associated  
23 segment pair, such as  $P(f|ph)$ ,  $P(i|y)$ , and so on. The error model probability  $P(f i s$   
24  $ik le | ph y s ic al)$  can then be expressed as the product of these segment pair  
25 probabilities, as follows:

$$P(\text{f i s i k l e} \mid \text{p h y s i c a l}) = P(\text{f} \mid \text{ph}) * P(\text{i} \mid \text{y}) * P(\text{s} \mid \text{s}) * P(\text{ik} \mid \text{ic}) * P(\text{le} \mid \text{al}).$$

The error model 44 examines all probabilities for all partitions over all possible words and selects the word that returns the highest probability, summed over all partitions. For example, let  $\text{Part}(w)$  be the set of all possible ways of partitioning word  $w$  and  $\text{Part}(s)$  be the set of all possible ways of partitioning the entered string  $s$ . For a particular partition  $R$  of the word  $w$  (i.e.,  $R \in \text{Part}(w)$ , where  $|R| = j$  contiguous segments), let partition  $R_i$  be the  $i^{\text{th}}$  segment. Similarly, for a particular partition  $T$  of the string  $s$  (i.e.,  $T \in \text{Part}(s)$ , where  $|T| = j$  contiguous segments), let partition  $T_i$  be the  $i^{\text{th}}$  segment. The error model 44 computes:

$$P(s \mid w) = \sum_{R \in \text{Part}(w)} P(R \mid w) \sum_{\substack{T \in \text{Part}(s) \\ |T|=|R|}} \prod_{i=1}^{|R|} P(T_i \mid R_i)$$

The first summation sums probabilities over all possible partitions of the word  $w$ . For a given partition  $R$ , the second summation sums over all possible partitions of the string  $s$ , with the restriction that both partitions must have the same number of segments. The product then multiplies the probabilities of each  $R_i \rightarrow T_i$ .

To demonstrate this computation, consider the word “physical” and the five-segment partition. The error model 44 tries different partitions  $R$  of the word, and for each word partition  $R$ , tries different partitions  $T$  of the entered string. For

each combination, the error model 44 computes a corresponding probability  $P(s|w)$ , as illustrated in Table 1:

Table 1

<u>Partitions</u>	<u>Probabilities <math>P(s w)</math></u>
$R_1$ : <u>ph</u> <u>y</u> <u>s</u> <u>ic</u> <u>al</u>	
$T_1$ : <u>f</u> <u>i</u> <u>s</u> <u>ik</u> <u>le</u>	$P(f ph) * P(i y) * P(s s) * P(ik ic) * P(le al)$
$T_2$ : <u>fi</u> <u>s</u> <u>i</u> <u>k</u> <u>le</u>	$P(fi ph) * P(s y) * P(i s) * P(k ic) * P(le al)$
$T_3$ : <u>fis</u> <u>i</u> <u>k</u> <u>le</u>	$P(fis ph) * P(i y) * P(l s) * P(l ic) * P(e al)$
$T_4$ : <u>fis</u> <u>ik</u> <u>le</u> _ _	$P(fis ph) * P(ik y) * P(le s) * P( ic) * P( al)$
$T_5$ : ...	....
$R_2$ : <u>phy</u> <u>si</u> <u>cal</u> _ _	
$T_1$ : <u>f</u> <u>i</u> <u>s</u> <u>ik</u> <u>le</u>	$P(f phy) * P(i si) * P(s cal) * P(ik  ) * P(le  )$
$T_2$ : <u>fi</u> <u>s</u> <u>i</u> <u>k</u> <u>le</u>	$P(fi phy) * P(s si) * P(i cal) * P(k  ) * P(le  )$
$T_3$ : <u>fis</u> <u>i</u> <u>k</u> <u>le</u>	$P(fis phy) * P(i si) * P(l cal) * P(l  ) * P(e  )$
$T_4$ : <u>fis</u> <u>ik</u> <u>le</u> _ _	$P(fis phy) * P(ik si) * P(le cal) * P(  ) * P(  )$
$T_5$ : ...	....
$R_2$ : ...	

After all permutations have been computed for a five-segment partition, the error model repeats the process for partitions of more or less than five segments. The error model 44 selects the word that yields the highest probability  $P(s|w)$ , summed over all possible partitions. The spell checker 40 uses the error model probability, along with the source model probability, to determine whether to autocorrect the entered string, leave the string alone, or suggest possible alternate words for the user to choose from.

1 If computational efficiency is a concern, the above relationship may be  
2 approximated as follows:

$$3 \\ 4 \quad P(s|w) = \max_{R \in \text{Part}(w), T \in \text{Part}(s)} P(R|w) * \prod_{i=1}^{|R|} P(T_i|R_i) \\ 5$$

6  
7 Two further simplifications can be made during implementation that still  
8 provides satisfactory results. One simplification is to drop the term  $P(R|w)$ .  
9 Another simplification is to set the terms  $P(T_i|R_i) = 1$  whenever  $T_i = R_i$ .

10 The error model 44 has a number of advantages over previous approaches.  
11 First, the error model is not constrained to single character edits, but robustly  
12 handles conversion of one arbitrary length string to another arbitrary length string.  
13 As noted in the Background, virtually all conventional spell checkers based on the  
14 noisy channel model use a fixed set of single character edits—insertion, deletion,  
15 and substitution—with some checkers also including simple transposition,  
16 doubling, and halving. However, people often mistype one string for another  
17 string, where one or both of the strings has length greater than one. These types of  
18 errors cannot be modeled succinctly using the conventional Weighted Levenshtein  
19 Distance.

20 The error model 44 captures the traditional single character edits,  
21 transposition, doubling, and halving, as well as many phenomena not captured in  
22 such simpler models. For example, if a person mistypes “philosophy” as  
23 “filosofy”, the error model 44 captures this directly by the edit “ph → f”, whereby  
24 a two-character string is converted to a single character string. Even when an  
25 error is a single letter substitution, often the environment in which it occurs is

1 significant. For instance, if a user enters “significant” as “significient”, it makes  
2 more sense to describe this by the edit operation “ant  $\rightarrow$  ent” than simply by “a  $\rightarrow$   
3 e”.

4 Another advantage of the error model 44 is that it can implement an even  
5 richer set of edits by allowing an edit to be conditioned on the position that the edit  
6 occurs,  $P(\alpha \rightarrow \beta \mid \text{PSN})$ , where PSN describes positional information about the  
7 substring within the word. For example, the position may be the start of a word,  
8 the end of a word, or some other location within the word (i.e.,  $\text{PSN} = \{\text{start of}$   
9  $\text{word, end of word, other}\}$ ). The spell checker adds a start-of-word symbol and an  
10 end-of-word symbol to each word to provide this positional information.

11 Fig. 2 shows the process implemented by the spell checker 40 to compute a  
12 probability  $P(s|w)$  given an entered string  $s$  and a dictionary word  $w$ . At block  
13 202, the spell checker receives a user-entered string  $s$  that might contain errors.  
14 Assume that the entered word is “fisikle”. Given this entered string  $s$ , the spell  
15 checker iterates over all words  $w$ . Suppose, for sake of discussion, the current  
16 word  $w$  is “physical”.

17 At block 204, the word  $w$  is partitioned into multiple segments. The word  
18 “physical” is partitioned, for example, into five segments “ph y s ic al”. At block  
19 206, the string  $s$  is partitioned into the same number of segments, such as “f i s ik  
20 le”.

21 At block 208, the error model computes a probability for this pair of  
22 partitioned strings as  $P(f|ph) * P(i|y) * P(s|s) * P(ik|ic) * P(le|al)$  and temporarily  
23 stores the result. The error model considers other partitions of the user-entered  
24 string  $s$  against the partitioned word  $w$ , as represented by the inner loop blocks 210  
25 and 212. With each completion of all possible partitions of the user-entered string

1 s, the error model 44 iteratively tries different partitions of the word  $w$ , as  
2 represented by the outer loop blocks 214 and 216.

3 At block 218, when all possible combinations of partitions have been  
4 processed, the error model 44 sums the probabilities to produce the probability  
5  $P(s|w)$ .

### 7 Training the Error Model

8 The error model 44 is trained prior to being implemented in the spell  
9 checker 40. The training is performed by a computing system, which may be the  
10 same computer as shown in Fig. 1 (i.e., the model is trained on the fly) or a  
11 separate computer employed by the developer of the spell checker (i.e., the model  
12 is trained during development). The training utilizes a training set or corpus that  
13 includes correct dictionary words along with errors observed when a user enters  
14 such words. One technique for training the error model is to use a training set  
15 consisting of <wrong, right> training pairs. Each training pair represents a  
16 spelling error together with the correct spelling of the word.

17 Fig. 3 shows a training computer 300 having a processor 302, a volatile  
18 memory 304, and a non-volatile memory 306. The training computer 300 runs a  
19 training program 308 to produce probabilities of different arbitrary-length string-  
20 to-string corrections ( $\alpha \rightarrow \beta$ ) over a large set of training words and associated  
21 mistakes observed from entry of such words. The training program 308 is  
22 illustrated as executing on the processor 302, although it is loaded into the  
23 processor from storage on non-volatile memory 306.

24 Training computer 300 has a training set 310 stored in non-volatile memory  
25 306 (i.e., hard disk(s), CD-ROM, etc.). The training set has <wrong, right>



1 training pairs. As an example, the training set 310 may have 10,000 pairs. The  
2 training computer uses the training set to derive probabilities associated with how  
3 likely the right word is to be changed to the wrong word. The probabilities are  
4 based on the least cost way to edit an arbitrary length character sequence  $\alpha$  into  
5 another arbitrary length character sequence  $\beta$  (i.e.,  $\alpha \rightarrow \beta$ ).

6 Fig. 4 shows a training method implemented by the training program 308.  
7 At block 402, the training method arranges the wrong and right words according to  
8 single letter edits: insertion (i.e.,  $\emptyset \rightarrow a$ ), substitution (i.e.,  $a \rightarrow b$ ), deletion (i.e.,  
9  $a \rightarrow \emptyset$ ), and match (i.e.,  $a = a$ ). The edits are assigned different weights. Using the  
10 Levenshtein Distance, for example, a match is assigned a weight of 0 and all other  
11 edits are given a weight of 1. Given a <Wrong, Right> training pair, the training  
12 method finds the least-cost alignment using single letter edits and edit weights.

13 Fig. 5 shows one possible least-cost alignment of a training pair <akgsual,  
14 actual>. In the illustrated alignment, the first letters "a" in each string match, as  
15 represented by the label "Mat" beneath the associated characters. This edit type is  
16 assessed a weight of 0. The second letters do not match. A substitution edit needs  
17 to be performed to convert the "c" into the "k", as represented by the legend  
18 "Sub". This edit type is assigned a weight of 1. There is no letter in the right  
19 word "actual" that corresponds to the letter "g" in the wrong word "akgsual" and  
20 hence an insertion edit is needed to insert the "g". Insertion is represented by the  
21 legend "Ins" and is given a weight of 1. Another substitution is needed to convert  
22 the "t" into an "s", and this substitution is also assessed a weight of 1. The last  
23 three letters in each string are matched and are accorded a weight of 0.

24 The alignment in Fig. 5 is one example of a least-cost alignment, having an  
25 edit cost of 3. Other alignments with the same cost may exist. For instance,

perhaps the letter "g" in "akgsual" may be aligned with "t" and "s" with space. This alternate alignment results in the same cost of 3. Selection of one alignment in such ties is handled arbitrarily.

After this initial alignment, all contiguous non-match edits are collapsed into a single error region (block 404 in Fig. 4). There may be multiple error regions in a given training pair, but the contiguous non-match edits are combined as common regions. Using the alignment of training pair <akgsual, actual> in Fig. 5, the contiguous "substitution-insertion-substitution" edits are collapsed into a single substitution edit "ct→kgs".

Fig. 6 shows the training pair <akgsual, actual> after all contiguous non-match edits are collapsed. Now, there is only one non-match edit, namely a generic substitution operation "ct→kgs".

An alternate way of training is to not collapse contiguous non-match edits. Given the alignment shown in Fig. 5, this would result in three substitution operations: c→k, NULL→g and t→s, instead of the single substitution operation obtained by collapsing.

To allow for richer contextual information, each substitution is expanded to incorporate one or more edits from the left and one or more edits from the right (block 406 in Fig. 4). As an example, the expansion might entail up to two edits from the left and two edits from the right. For the substitution "ct→kgs", the training method generates the following substitutions:

ct→kgs

act→akgs

actu→akgsu

1                   ctu→kgsu

2                   ctua→kgsua

3  
4           Each of these possible substitutions is assigned an equal fractional count,  
5 such as one-fifth of a count per substitution.

6           At block 408, the probability of each substitution  $\alpha \rightarrow \beta$  is computed as  
7  $\text{count}(\alpha \rightarrow \beta) / \text{count}(\alpha)$ . For instance, to compute  $P(\text{ct} \rightarrow \text{kgs})$ , the method first  
8 sums up all of the counts found for the edit  $\text{ct} \rightarrow \text{kgs}$  in the training corpus. Then,  
9 the method counts the number of times the substring "ct" is seen in a suitably  
10 sized corpus of representative text, and divides the first count by the second.

11           After obtaining the set of edits and edit probabilities, the process may  
12 iteratively re-estimate the probabilities using a form of the well known E-M  
13 algorithm, similar to the retraining method described in the Church paper.  
14 However, the inventors have observed that very good results can be obtained  
15 without re-estimating the parameters.

16           By varying the training set 310, the error model 44 may be trained to  
17 accommodate the error profiles on a user-by-user basis, or on a group-by-group  
18 basis. For instance, a user with dyslexia is likely to have a very different error  
19 profile than somebody without dyslexia. An English professor is likely to have a  
20 very different error profile from a third grader, and a native Japanese speaker  
21 entering English text is likely to have a very different error profile from a native  
22 English speaker.

23           Therefore, the efficacy of the spelling correction program can be improved  
24 further if it is trained to the particular error profile of an individual or  
25 subpopulation. For a relatively static subpopulation, the training set 310 is created

1 to contain <wrong, right> pairs from the subpopulation. The error model 44 is  
2 then trained based on this training set.

3 For individuals, a generally trained error model can be configured to adapt  
4 to the user's own tendencies. As the user employs the spell checker, it keeps track  
5 of instances where an error is corrected. One way to track such instances is to  
6 monitor which word the user accepts from a list of corrections presented by the  
7 spell checker when it flags a word as incorrect. Another way is to monitor when  
8 the spell checker autocorrects a string the user has input. By tracking corrected  
9 errors, the spell checker collects <wrong, right> pairs that are specific to that  
10 individual. This can then be used to adapt the error model to the individual, by  
11 retraining the  $(\alpha \rightarrow \beta)$  parameters to take into account these individual error tuples.

12 It is desirable to use a large number of <wrong, right> pairs for training, as  
13 this typically improves accuracy of the resultant correction probabilities. One  
14 method for collecting the training pairs is to harvest it from available on-line  
15 resources such as the World Wide Web. A spell checker can auto-correct a string  
16  $s$  into a word  $w$  when it is sufficiently certain that  $w$  was the intended word that  
17 was mistyped as  $s$ . For instance, the spell checker can auto-correct  $s$  into  $w$  if  $w$  is  
18 the most likely intended word according to our model, and the second most likely  
19 intended word is sufficiently less probable than  $w$ . The error model can thus be  
20 iteratively trained as follows:

- 21  
22 (1) Obtain a set of <wrong, right> pairs and use them to train an initial  
23 model.
- 24 (2) Run the model over a collection of on-line resources. In all cases, when  
25 the model auto-corrects string  $s$  into word  $w$ , save the tuple  $\langle s, w \rangle$ .

1 (3) Use these saved  $\langle s, w \rangle$  tuples to retrain the model.

2 (4) Go to step (2).

### 3 4 Spell Correction Method

5 Once the error model is trained, the spell checker 40 is ready to identify and  
6 correct misspelled words. As noted earlier, the spell checker 40 attempts to  
7 correct an erroneously entered string  $s$  by returning the change that maximizes the  
8 probabilities of the source and error models, as follows:

$$\arg \max_{w \in D} P(s | w) \times P(w | context)$$

12 One approach to performing this search is to first return the  $k$  best candidate  
13 words according to  $P(s|w)$  and then re-score these  $k$  best words according to the  
14 full model  $P(s|w) * P(w|context)$ .

15 To find  $\arg \max_{w \in D} P(s | w)$ , the spell checker can be configured to iterate  
16 over the entire dictionary  $D$ . However, it is much more efficient to convert the  
17 dictionary into a trie and compute edit costs for each node in the trie.  
18 Representing a dictionary as a trie is conventional and well known to those of skill  
19 in the art. Further efficiency gains can be had if the set of edits are stored as a trie  
20 of edit left-hand-sides, with pointers to corresponding tries of right-hand-sides of  
21 edit rules.

22 Depending upon the certainty that word  $w$  is intended when string  $s$  is  
23 input, the spell checker 40 has the following options: (1) leave the string  
24 unchanged, (2) autocorrect the string  $s$  into the word  $w$ , or (3) offer a list of  
25

1 possible corrections for the string *s*. The spell correction process may be  
2 represented by the following pseudo code:

3       For each space-delimited string *s*

4             Find the *k* most likely words

5             If there is sufficient evidence that *s* is the intended string

6                 Then do nothing

7             Else

8                 If there is sufficient evidence that the most likely word *w*  
9                 is the intended word given the generated string *s*,

10                 Then autocorrect *s* into *w*

11             Else

12                 Flag the string *s* as potentially incorrect and offer  
13                 the user a sorted list of possible corrections.

#### 14       Conclusion

15       Although the description above uses language that is specific to structural  
16 features and/or methodological acts, it is to be understood that the invention  
17 defined in the appended claims is not limited to the specific features or acts  
18 described. Rather, the specific features and acts are disclosed as exemplary forms  
19 of implementing the invention.  
20  
21  
22  
23  
24  
25

1 **CLAIMS**

2 We claim:

3  
4 **1.** A method comprising:

5 receiving an entered string; and

6 determining how likely a word was to have been entered as the string based  
7 on at least one edit operation that converts a first character sequence of arbitrary  
8 length in the word to a second character sequence of arbitrary length in the string.  
9

10 **2.** A method as recited in claim 1, wherein the first character sequence  
11 has a first length and the second character sequence has a second length that is  
12 different than the first length.  
13

14 **3.** A method as recited in claim 1, wherein the first character sequence  
15 has multiple characters and the second character sequence has multiple characters.  
16

17 **4.** A method as recited in claim 1, wherein the first character sequence  
18 has a first number of multiple characters and the second character sequence has a  
19 second number of multiple characters that is different from the first number of  
20 multiple characters.  
21

22 **5.** A method as recited in claim 1 and further comprising determining  
23 how likely the word is to have been generated.  
24  
25

1           6.    A method as recited in claim 1 and further comprising conditioning  
2 the edit operation on a position that the edit occurs at within the word.

3  
4           7.    A method as recited in claim 1 and further comprising identifying the  
5 string as potentially incorrect.

6  
7           8.    A method as recited in claim 1 and further comprising correcting the  
8 string to the word.

9  
10          9.    A computer readable medium having computer-executable  
11 instructions that, when executed on a processor, perform the method as recited in  
12 claim 1.

13  
14          10.   A method comprising:  
15       receiving an entered string  $s$ ; and  
16       determining a probability  $P(s|w)$  expressing how likely a word  $w$  was to  
17 have been incorrectly entered as the string  $s$  based on one or more edit operations  
18 that convert first arbitrary-length character sequences  $\alpha_1, \alpha_2, \alpha_3, \dots, \alpha_n$  in the  
19 word  $w$  to corresponding second arbitrary-length character sequences  $\beta_1, \beta_2, \beta_3,$   
20  $\dots, \beta_n$  in the string  $s$ , wherein:

21  
22                   
$$P(s|w) = P(\beta_1|\alpha_1) * P(\beta_2|\alpha_2) * P(\beta_3|\alpha_3) * \dots * P(\beta_n|\alpha_n)$$
  
23  
24  
25



1           **11.**    A method as recited in claim 10, wherein lengths of corresponding  
2 first and second character sequences are different.

3  
4           **12.**    A method as recited in claim 10 and further comprising determining  
5 how likely the word  $w$  is to have been generated.

6  
7           **13.**    A method as recited in claim 10 and further comprising conditioning  
8 the edit operations on positions that the edits occur at within the word.

9  
10          **14.**    A method as recited in claim 10 and further comprising correcting  
11 the string  $s$  to the word  $w$ .

12  
13          **15.**    A method as recited in claim 10 and further comprising identifying  
14 the string  $s$  as potentially incorrect.

15  
16          **16.**    A computer readable medium having computer-executable  
17 instructions that, when executed on a processor, perform the method as recited in  
18 claim 10.

19  
20          **17.**    A method comprising:  
21 receiving an entered string  $s$ ; and  
22 determining a probability  $P(s|w)$  expressing how likely a word  $w$  was to  
23 have been incorrectly entered as the string  $s$ , by partitioning the word  $w$  and the  
24 string  $s$  and computing probabilities for various partitionings, as follows:  
25

$$P(s | w) = \sum_{R \in \text{Part}(w)} P(R | w) \sum_{\substack{T \in \text{Part}(s) \\ |T|=|R|}} \prod_{i=1}^{|R|} P(T_i | R_i)$$

where  $\text{Part}(w)$  is a set of possible ways of partitioning the word  $w$ ,  $\text{Part}(s)$  is a set of possible ways of partitioning the string  $s$ ,  $R$  is a particular partition of the word  $w$ , and  $T$  is a particular partition of the string  $s$ .

18. A method as recited in claim 17 and further comprising selecting the partition that returns a highest probability.

19. A method as recited in claim 17 and further comprising determining how likely the word  $w$  is to have been generated.

20. A method as recited in claim 17 and further comprising correcting the string  $s$  to the word  $w$ .

21. A method as recited in claim 17 and further comprising identifying the string  $s$  as potentially incorrect.

22. A computer readable medium having computer-executable instructions that, when executed on a processor, perform the method as recited in claim 17.

1       **23.**    A method comprising:

2       receiving an entered string  $s$ ; and

3       determining a probability  $P(s|w)$  expressing how likely a word  $w$  was to  
4       have been incorrectly entered as the string  $s$ , by partitioning the word  $w$  and the  
5       string  $s$  and computing probabilities for various partitionings, as follows:

6

$$7 \quad P(s|w) = \max_{R \in \text{Part}(w), T \in \text{Part}(s)} P(R|w) * \prod_{i=1}^{|R|} P(T_i|R_i)$$

8

9

10       where  $\text{Part}(w)$  is a set of possible ways of partitioning the word  $w$ ,  $\text{Part}(s)$  is a set  
11       of possible ways of partitioning the string  $s$ ,  $R$  is a particular partition of the word  
12        $w$ , and  $T$  is a particular partition of the string  $s$ .

13

14       **24.**    A method as recited in claim 23 and further comprising omitting the  
15       term  $P(R|w)$  from the computation of  $P(s|w)$ .

16

17       **25.**    A method as recited in claim 23 and further comprising setting  
18       terms  $P(T_i|R_i) = 1$  whenever  $T_i = R_i$ .

19

20       **26.**    A method as recited in claim 23 and further comprising determining  
21       how likely the word  $w$  is to have been generated.

22

23       **27.**    A method as recited in claim 23 and further comprising correcting  
24       the string  $s$  to the word  $w$ .

25

1           **28.**    A method as recited in claim 23 and further comprising identifying  
2 the string  $s$  as potentially incorrect.

3  
4           **29.**    A computer readable medium having computer-executable  
5 instructions that, when executed on a processor, perform the method as recited in  
6 claim 23.

7  
8           **30.**    A method comprising:  
9           receiving an entered string  $s$ ; and  
10          determining a probability  $P(s|w)$  expressing how likely a word  $w$  was to  
11 have been incorrectly entered as the string  $s$ , by partitioning the word  $w$  and the  
12 string  $s$  and finding a partition  $R$  of the word  $w$  and a partition  $T$  of the string  $s$   
13 such that  $\prod_{i=1}^{|R|} P(T_i | R_i)$  is maximized.

14  
15          **31.**    A method as recited in claim 30 and further comprising determining  
16 how likely the word  $w$  is to have been generated.

17  
18          **32.**    A method as recited in claim 30 and further comprising correcting  
19 the string  $s$  to the word  $w$ .

20  
21          **33.**    A method as recited in claim 30 and further comprising identifying  
22 the string  $s$  as potentially incorrect.  
23  
24  
25

1       **34.**   A computer readable medium having computer-executable  
2 instructions that, when executed on a processor, perform the method as recited in  
3 claim 30.

4  
5       **35.**   A method for training an error model used in a spell checker,  
6 comprising:

7       determining, given a <wrong, right> training pair and multiple single  
8 character edits that convert characters in one of the right or wrong strings to  
9 characters in the other of the right or wrong strings at differing costs, an alignment  
10 of the wrong string and the right string that results is a least cost to convert the  
11 characters;

12       collapsing any contiguous non-match edits into one or more common error  
13 regions, each error region containing one or more characters that can be converted  
14 to one or more other characters using a substitution edit; and

15       computing a probability for each substitution edit.

16  
17       **36.**   A method as recited in claim 35, wherein the assigning comprises  
18 assessing a cost of 0 to all match edits and a cost of 1 to all non-match edits.

19  
20       **37.**   A method as recited in claim 35, wherein the single character edits  
21 comprises insertion, deletion, and substitution.

22  
23       **38.**   A method as recited in claim 35, further comprising collecting  
24 multiple <wrong, right> training pairs from online resources.  
25

1           **39.**    A method as recited in claim 35, further comprising expanding each  
2 of the error regions to capture at least one character on at least one side of the error  
3 region.

4  
5           **40.**    A program embodied on a computer readable medium, which when  
6 executed, directs a computer to perform the following:

7                receive an entered string; and

8                determine how likely an expected string was to have been entered as the  
9 entered string based on at least one edit operation that converts a first character  
10 sequence of arbitrary length in the expected string to a second character sequence  
11 of arbitrary length in the entered string.

12  
13           **41.**    A program as recited in claim 40, wherein the first character  
14 sequence has a first length and the second character sequence has a second length  
15 that is different than the first length.

16  
17           **42.**    A program as recited in claim 40, wherein the first character  
18 sequence has multiple characters and the second character sequence has multiple  
19 characters.

20  
21           **43.**    A program as recited in claim 40, wherein the first character  
22 sequence has a first number of multiple characters and the second character  
23 sequence has a second number of multiple characters that is different from the first  
24 number of multiple characters.  
25

1       **44.**   A program as recited in claim 40, further comprising computer-  
2 executable instructions that directs a computer to determine how likely the  
3 expected string is to have been generated.

4  
5       **45.**   A program as recited in claim 40, further comprising computer-  
6 executable instructions that directs a computer to perform, depending upon how  
7 likely an expected string was to be incorrectly entered as the entered string, one of  
8 the following: (1) leave the entered string unchanged, (2) autocorrect the entered  
9 string into the expected string, or (3) offer a list of possible corrections.

10  
11       **46.**   A spell checker program, embodied on a computer-readable  
12 medium, comprising the program of claim 40.

13  
14       **47.**   A language conversion program, embodied on a computer-readable  
15 medium, comprising the program of claim 40.

16  
17       **48.**   A word processing program, embodied on a computer-readable  
18 medium, comprising the program of claim 40.

19  
20       **49.**   A program embodied on a computer readable medium, which when  
21 executed, directs a computer to perform the following:

22       (1) receive an entered string  $s$ ;

23       (2) for multiple words  $w$  in a dictionary, determine:

24               (a) how likely a word  $w$  in a dictionary is to have been generated,

25                $P(w|context)$ ; and

1 (b) how likely the word  $w$  was to have been entered as the string  
2  $s$ ,  $P(s|w)$ , based on at least one edit operation that converts a first  
3 character sequence of arbitrary length in the word to a second  
4 character sequence of arbitrary length in the string; and

5 (3) maximize  $P(s|w)*P(w|context)$  to identify which of the words is most  
6 likely the word intended when the string  $s$  was entered.

7  
8 **50.** A program as recited in claim 49, wherein the determination (2) is  
9 performed for all words in the dictionary.

10  
11 **51.** A program as recited in claim 49, further comprising computer-  
12 executable instructions that directs a computer to either (1) leave the string  
13 unchanged, (2) autocorrect the string into the word, or (3) offer a list of possible  
14 corrections.

15  
16 **52.** A spell checker program, embodied on a computer-readable  
17 medium, comprising the program of claim 49.

18  
19 **53.** A language conversion program, embodied on a computer-readable  
20 medium, comprising the program of claim 49.

21  
22 **54.** A spell checker comprising:  
23 a source model component to determine how likely a word  $w$  in a  
24 dictionary is to have been generated; and  
25



1 an error model component to determine how likely the word  $w$  was to have  
2 been incorrectly entered as the string  $s$  based on arbitrary length string-to-string  
3 transformations.

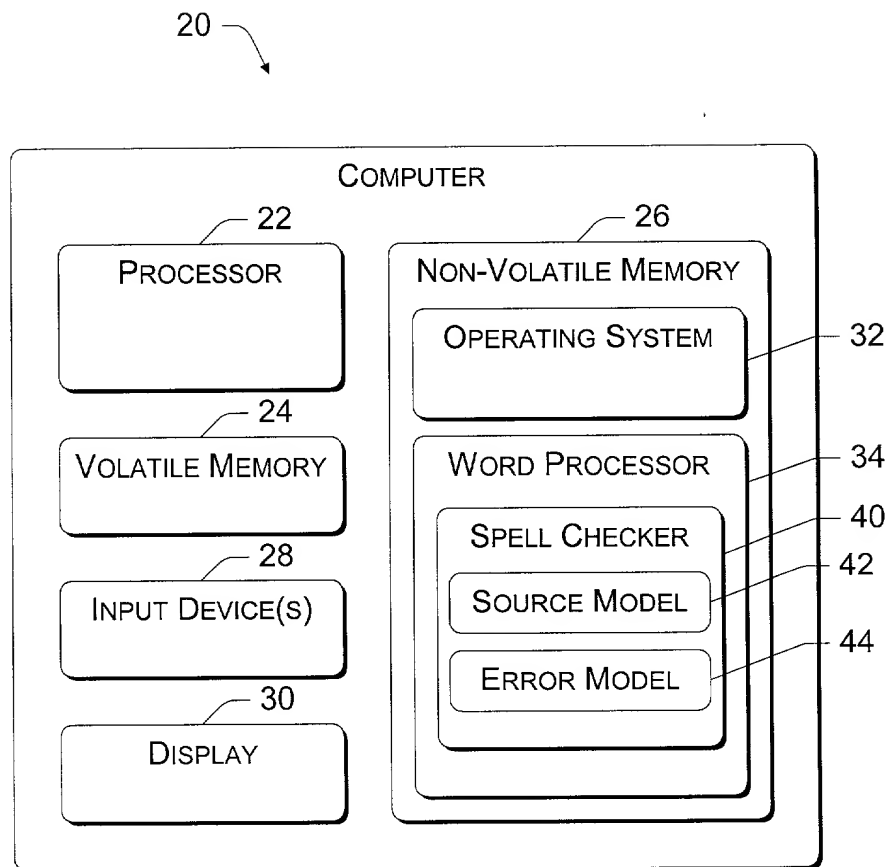
4  
5 **55.** A spell checker as recited in claim 54, wherein the string-to-string  
6 transformations involve conversion of a first character sequence of a first length  
7 into a second character sequence of a second length that is different than the first  
8 length.

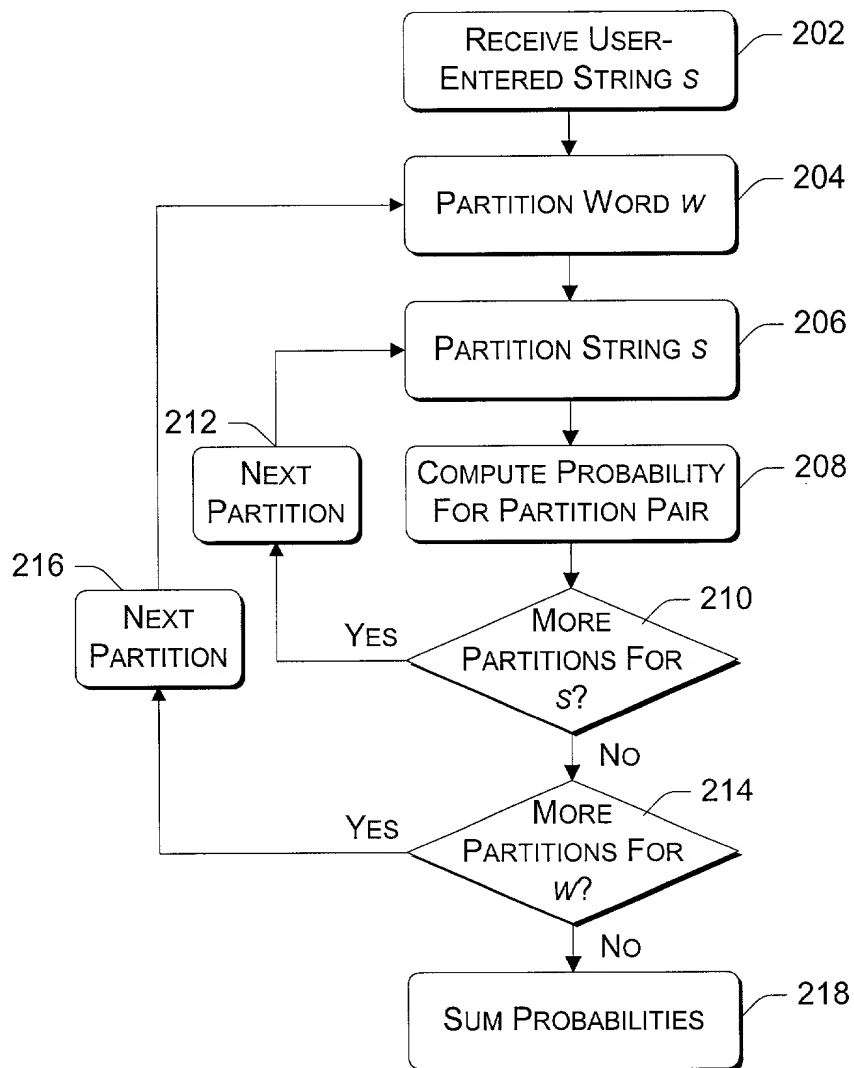
9  
10 **56.** A spell checker as recited in claim 54, wherein the string-to-string  
11 transformations involve conversion of a first character sequence with multiple  
12 characters into a second character sequence with multiple characters.

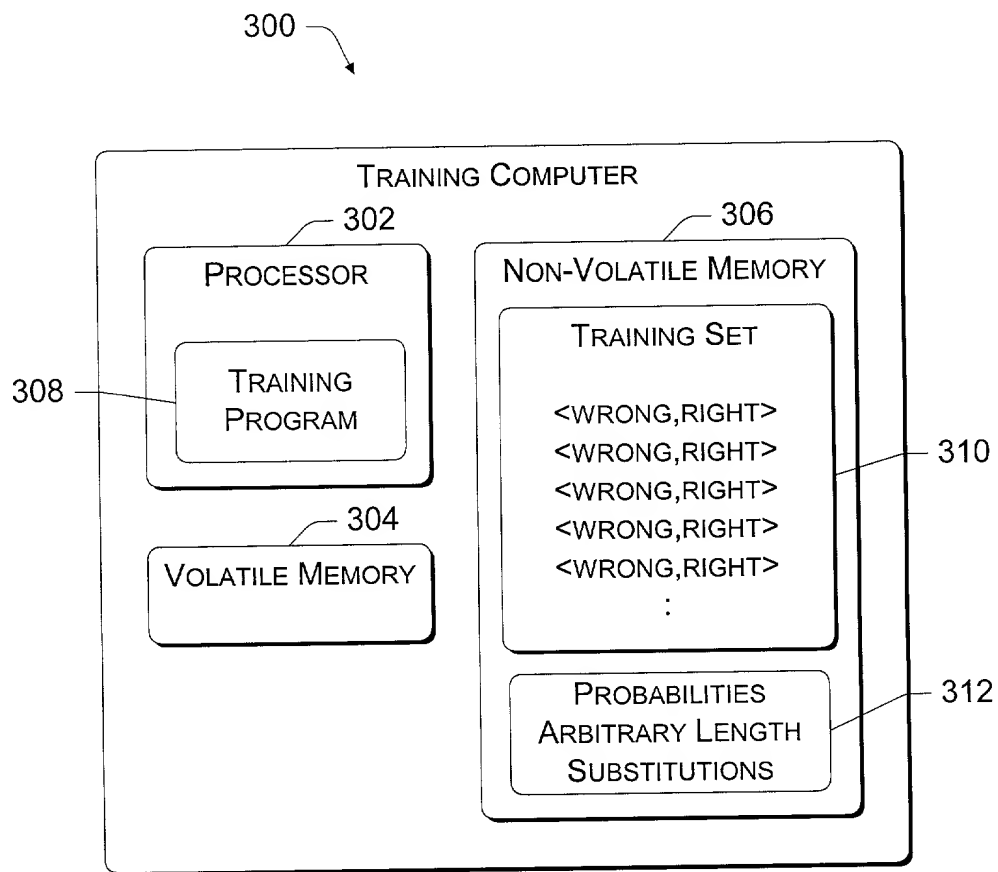
13  
14 **57.** A spell checker as recited in claim 54, wherein the string-to-string  
15 transformations involve conversion of a first character sequence having a first  
16 number of multiple characters into a second character sequence having a second  
17 number of multiple characters that is different from the first number of multiple  
18 characters.

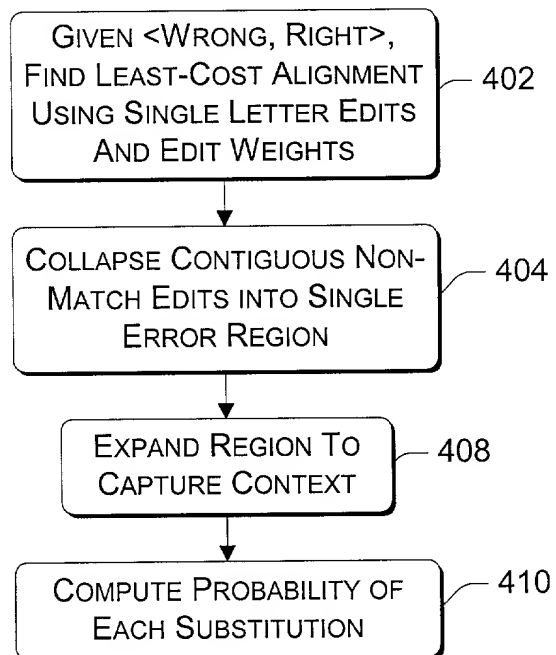
1 ABSTRACT

2 A spell checker based on the noisy channel model has a source model and  
3 an error model. The source model determines how likely a word  $w$  in a dictionary  
4 is to have been generated. The error model determines how likely the word  $w$  was  
5 to have been incorrectly entered as the string  $s$  (e.g., mistyped or incorrectly  
6 interpreted by a speech recognition system) according to the probabilities of  
7 string-to-string edits. The string-to-string edits allow conversion of one arbitrary  
8 length character sequence to another arbitrary length character sequence.  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25

*Fig. 1*

*Fig. 2*

*Fig. 3*

*Fig. 4*

	a	c		t	u	a	l
	↓	↓	↓	↓	↓	↓	↓
	a	k	g	s	u	a	l
<b>Edit</b>	Mat	Sub	Ins	Sub	Mat	Mat	Mat
<b>Weight</b>	0	1	1	1	0	0	0

*Fig. 5*

	a		ct		u	a	l
	↓		↓		↓	↓	↓
	a		kgs		u	a	l
<b>Edit</b>	Mat		Sub		Mat	Mat	Mat

*Fig. 6*